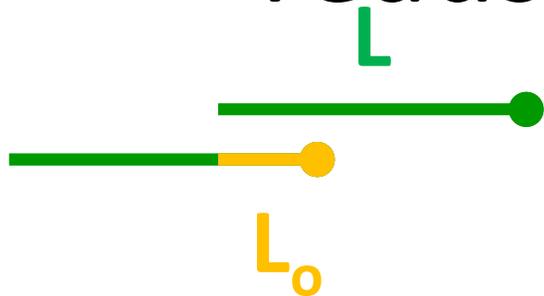


How long should be the length L_{ov} of the overlap to connect two short reads into a contig?



G

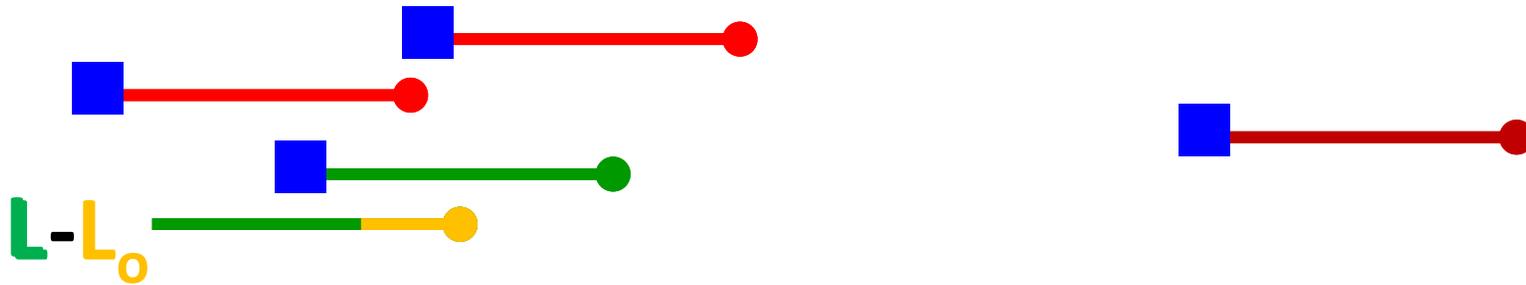
If DNA was a random chain with $p_A = p_C = p_G = p_T = 1/4$

$L_{ov} \sim 16-20$ would be enough

$$2 \cdot G \cdot 4^{-L_{ov}} = 2 \cdot 3 \times 10^9 \cdot 4^{-16} = 1.4$$

$$2 \cdot 3 \times 10^9 \cdot 4^{-20} = 0.0055 \ll 1$$

How many contigs?

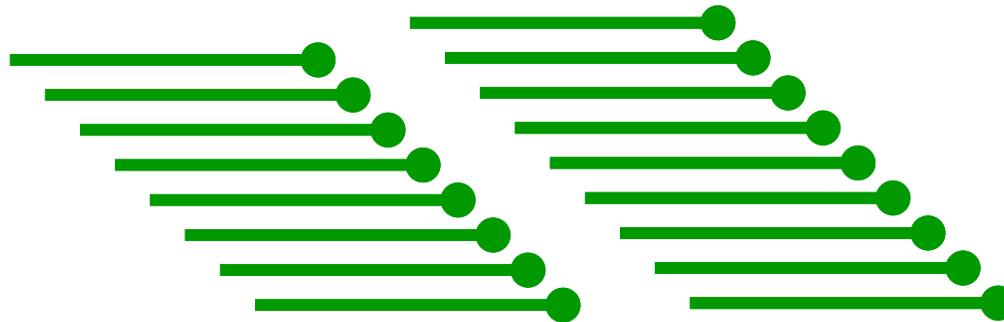


G

$$P(\text{short read can be extended by another short read}) = \frac{L - L_0}{G} = p$$

$$P(\text{short read cannot be extended by any short reads}) = e^{-pN} \approx Ne^{-\lambda}$$

$$\text{number of contigs} = Ne^{-pN} \approx Ne^{-\lambda}$$



How many contigs?

- A given short read is the right end of a contig if and only if no left ends of other short reads fall within it.
- The left end of another short read has the probability $p=(L-1)/G$ to fall within a given read. There are $N-1$ other reads. Hence the expected number of left ends inside a given shot read is $p \cdot (N-1)=(N-1) \cdot (L-1)/G \approx \lambda$
- If significant overlap required to merge two short reads is L_{ov} , modified λ is given by $(N-1) \cdot (L - L_{ov})/G$
- Probability that no left ends fall inside a short read is $exp(-\lambda)$. Thus the Number of contigs is $N_{contigs}=Ne^{-\lambda}$:

λ	0.5	0.75	1	1.5	2	3	4	5	6	7
Mean number of contigs	60.7	70.8	73.6	66.9	54.1	29.9	14.7	6.7	3.0	1.3

Table 5.2. The mean number of contigs for different levels of coverage, with $G = 100,000$ and $L = 500$.

Average length of a contig?

- Length of a genome covered:

$$G_{covered} = G \cdot P(X > 0) = G \cdot (1 - \exp(-\lambda))$$

- Number of contigs $N_{contigs} = N \cdot e^{-\lambda}$

- Average length of a contig =

$$\langle L \rangle = \sum_i L_i / N_{contigs} = G_{covered} / N_{contigs} =$$

$$G \cdot (1 - \exp(-\lambda)) / N \cdot e^{-\lambda} = L \cdot (1 - \exp(-\lambda)) / \lambda \cdot e^{-\lambda}$$

λ	2	4	6	8	10
Mean contig size	1,600	6,700	33,500	186,000	1,100,000

Table 5.3. The mean contig size for different values of a for the case $L = 500$.

Estimate

- Human genome is 3×10^9 bp long
- Chromosome 1 is about $G = 0.25 \times 10^9$ bp
- Illumina generates short reads $L = 100$ bp long
- What number of reads N are needed to completely assemble the 1st chromosome?
- The formula to use is: $1 = N_{contigs} = N e^{-\lambda} = N e^{-NL/G}$
- Answer: $N = 4.4 \times 10^7$ short (100bp) reads
Test: $4.4e7 * \exp(-4.4e7 * 100 / 0.25e9) = 0.99997$
- What coverage redundancy λ will it be?
Answer: $\lambda = NL/G = 17.6$ coverage redundancy

How much would it cost to assemble the human genome now?

- Human Genome Project: **\$2.7 billion** in 1991 dollars.
- Now a **de novo full assembly** of the whole human genome would now cost $3 \times 10^9 \times 17.6 / 10^9 \times 10\$/\text{GBase} = \$ 530$
- **2nd genome** (and after) would be **even cheaper** as we would already have a **reference genome to** which we can **map short reads**. (Puzzle: picture on the box)
- But this is a **naïve estimate**. In reality, there are complications. See the next slides:

What spoils these estimates?

```
>gi|224514922|ref|NT_024477.14| Homo sapiens chromosome 12 genomic
contig, GRCh37.p13 Primary Assembly (displaying 3' end)
CGGGAAATCAAAAGCCCCTCTGAATCCTGCGCACCGAGATTCTCCCCAGCCAAGGTGAGGCGGCAGCAGT
GGGAGATCCACACCGTAGCATTGGAACACAAATGCAGCATTACAAATGCAGACATGACACCGAAAATATA
ACACACCCCATTGCTCATGTAACAAGCACCTGTAATGCTAATGCACTGCCTCAAAACAAAATATTAATAT
AAGATCGGCAATCCGCACACTGCCGTGCAGTGCTAAGACAGCAATGAAAATAGTCAACATAATAACCCTA
ATAGTGTTAGGGTTAGGGTCAGGGTCCCGGTCCGGGTCCGGGTCCGGGTCCGGGTCCGGGTCCGGGTCA
GGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGT
TAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGG
GTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTA
GGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGT
TAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGG
GTTAGGGTTAGGGTTAGGGTTAG
```

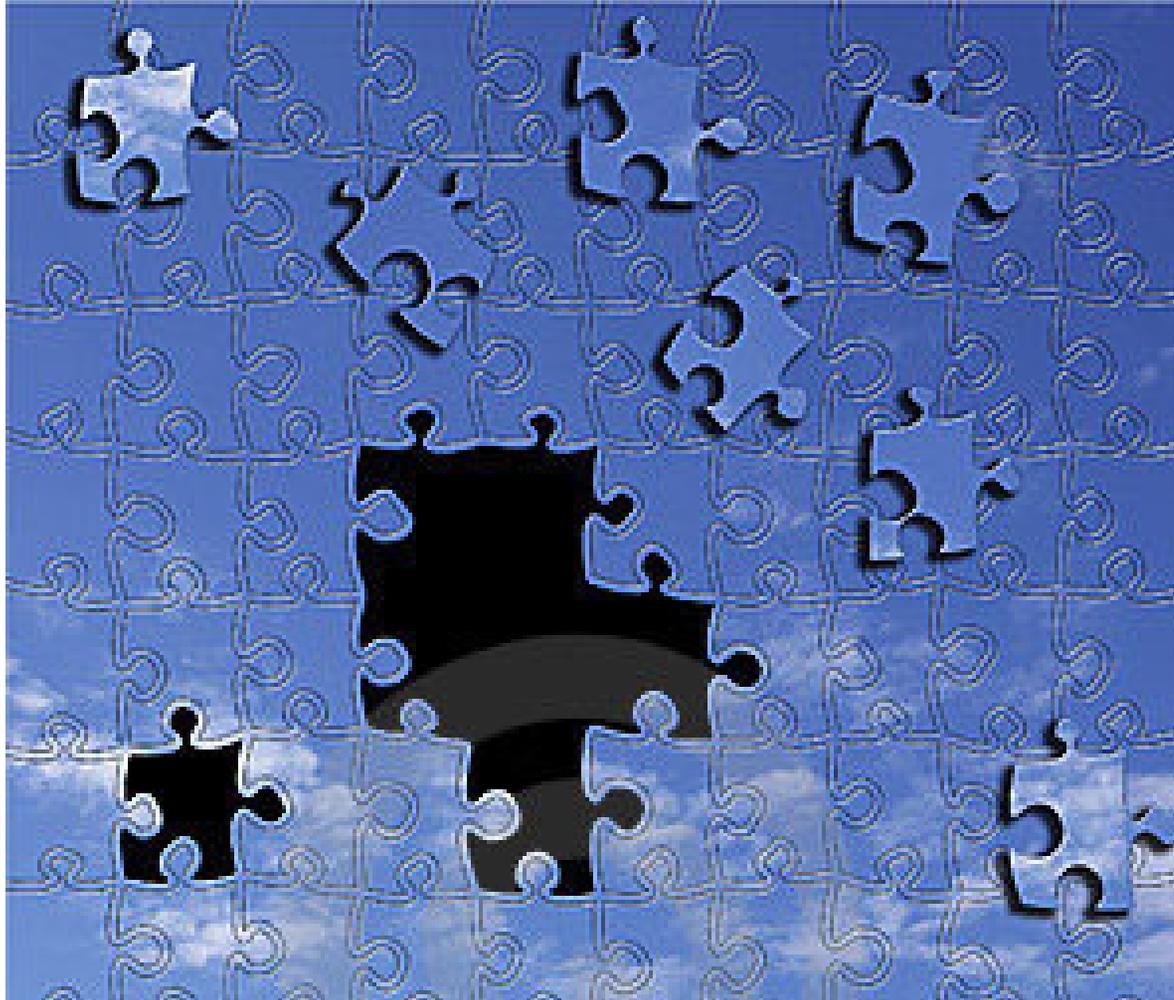
FIGURE 8.11 A BLASTN search of the human genome (all assemblies) database was performed at the NCBI website using **TTAGGGTTAGGGTTAGGG** as query (i.e., three TTAGGG repeats). There were matches to hundreds of genomic scaffolds. This figure shows an example (NT_024477.14) assigned to the **telomere of chromosome 12q having many dozens of TTAGGG repeats.** These occurred at the 3' end of the genomic contig sequence.

There were **100s of matches** while **one expects $\ll 1$ match:**

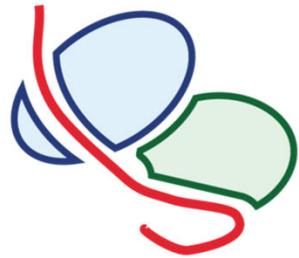
$$2 \cdot 3 \times 10^9 \cdot 4^{-18} = 0.08 \ll 1$$

DNA repeats make assembly difficult

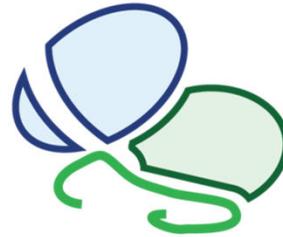
Repeats are like sky puzzle pieces



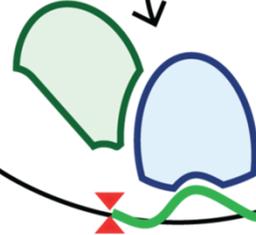
Formation of
Ribonucleoprotein complexes



Reverse
Transcription



Integration

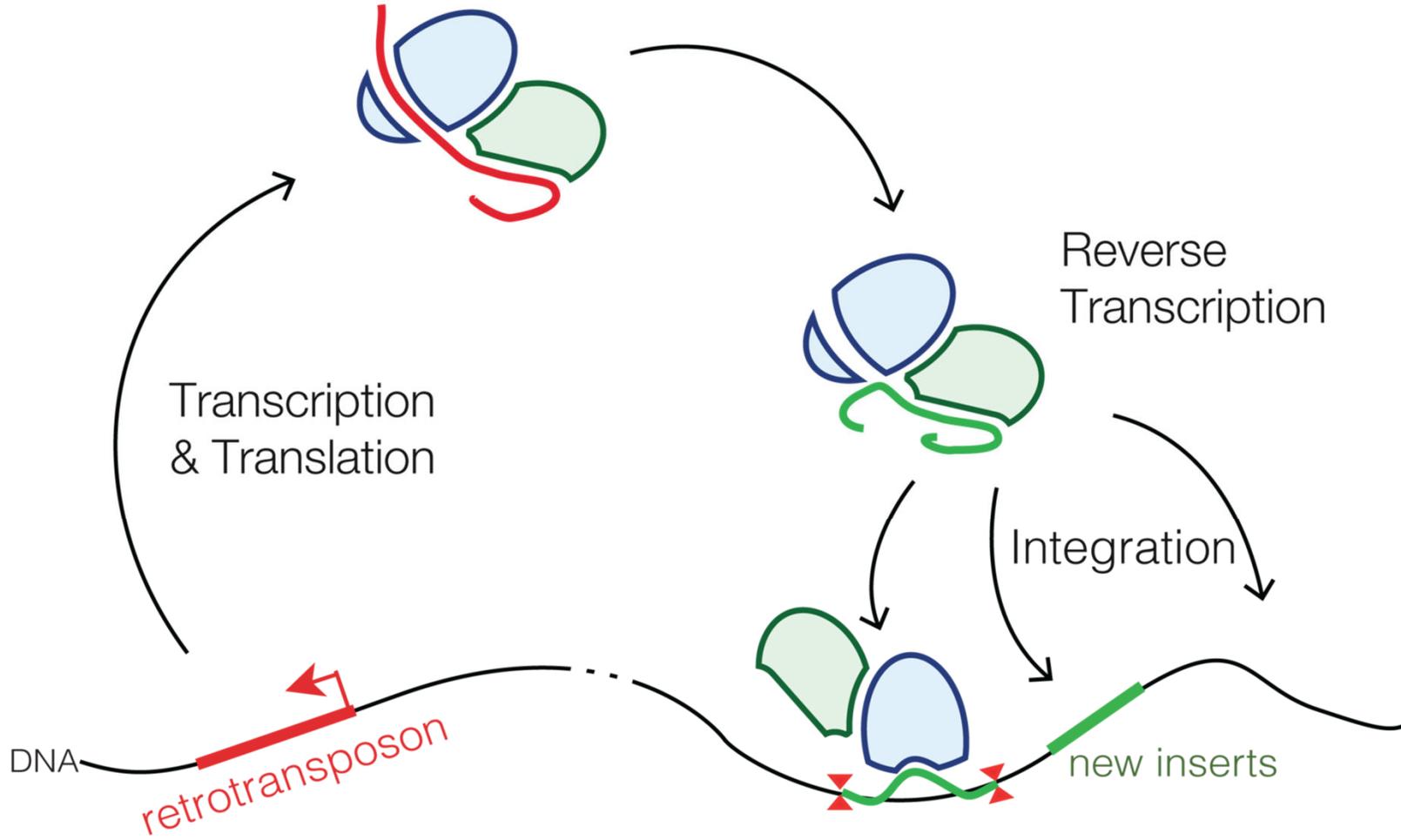


new inserts

Transcription
& Translation

DNA

retrotransposon



For each of three classes of **parasitic (selfish) DNA** there is a **shorter version** parasitizing the **longer version**

Classes of interspersed repeat in the human genome

			Length	Copy number	Fraction of genome	
Long Interspersed Nuclear Elements	LINEs	Autonomous		6–8 kb	850,000	21%
	SINEs	Non-autonomous		100–300 bp		
Retrovirus-like elements	Autonomous		6–11 kb	450,000	8%	
	Non-autonomous		1.5–3 kb			
DNA transposon fossils	Autonomous		2–3 kb	300,000	3%	
	Non-autonomous		80–3,000 bp			

Adapted from slide by Ross Hardison, Penn State U.

How to assemble a real genome with repeats?

Here we assume a “de novo” assembly
without help from the previously
assembled genomes



Nicolaas Govert de Bruijn (1918 – 2012) was a Dutch mathematician, noted for his many contributions in the fields of **graph theory**, analysis, number theory, combinatorics and logic

Courtesy of [Ben Langmead](#). Used with permission.

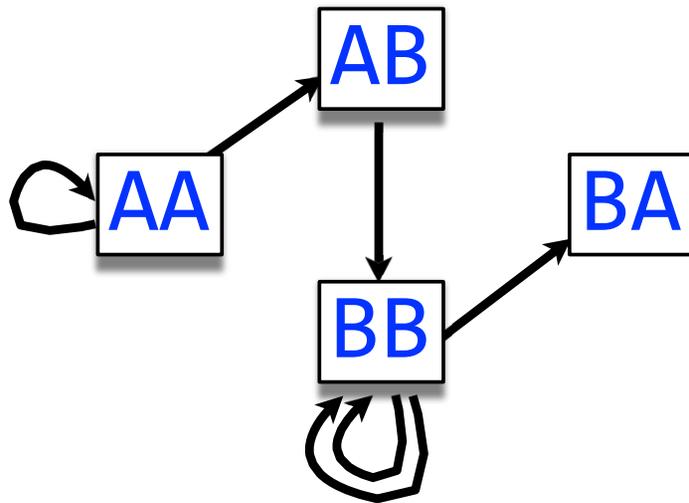
<http://www.langmead-lab.org/teaching-materials/>

De Bruijn graph

genome: **AAABBBBA**

3-mers: **AAA, AAB, ABB, BBB, BBB, BBA**

L/R 2-mers: **AA, AA AA, AB AB, BB BB, BB BB, BB BB, BA**



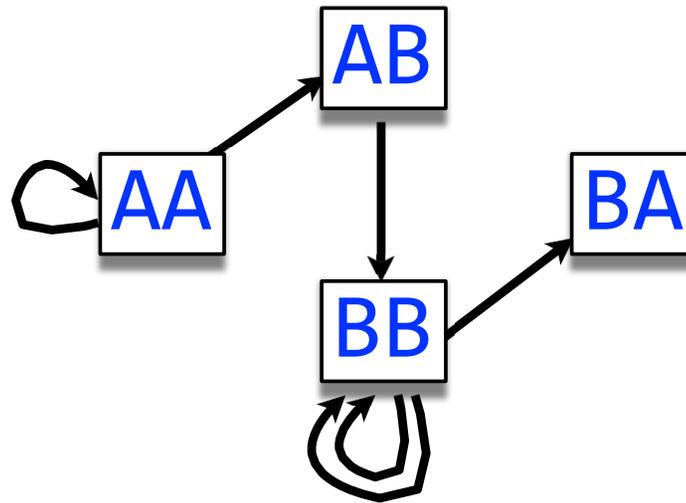
One edge per **every** k -mer

One node per **distinct** $k-1$ -mer

Courtesy of [Ben Langmead](#). Used with permission.

<http://www.langmead-lab.org/teaching-materials/>

De Bruijn graph

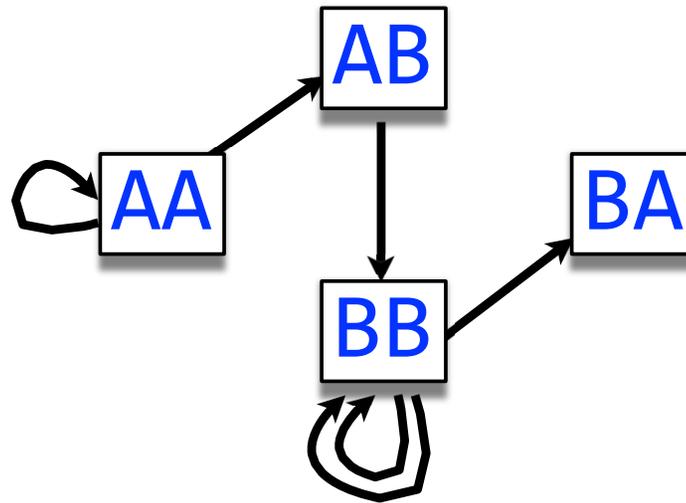


Walk crossing each edge exactly once gives a reconstruction of the genome

Courtesy of [Ben Langmead](#). Used with permission.

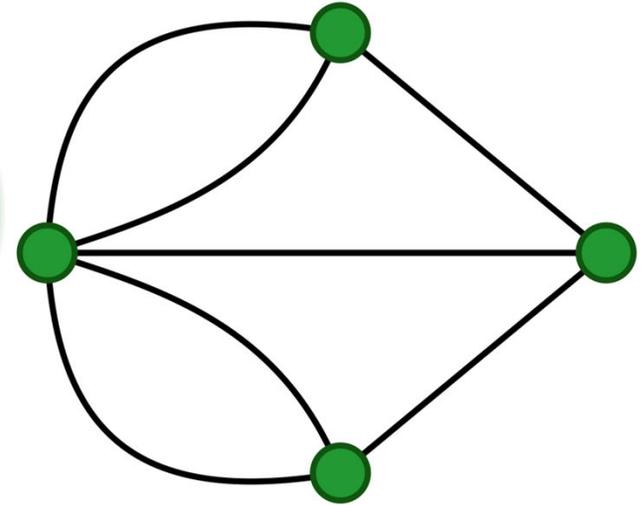
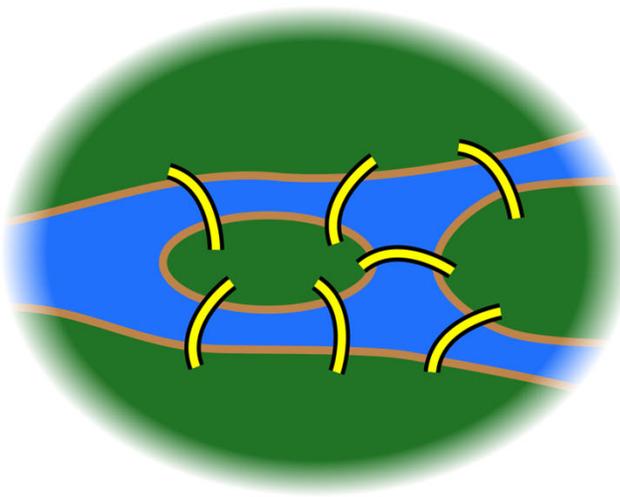
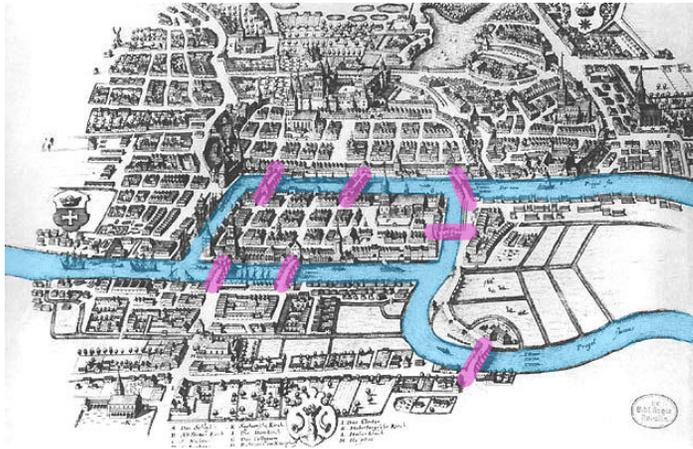
<http://www.langmead-lab.org/teaching-materials/>

Assembly = Eulerian walk on De Bruijn graph



AAABBBBA

Walk crossing each edge exactly once gives a reconstruction of the genome. This is a *Eulerian walk*.

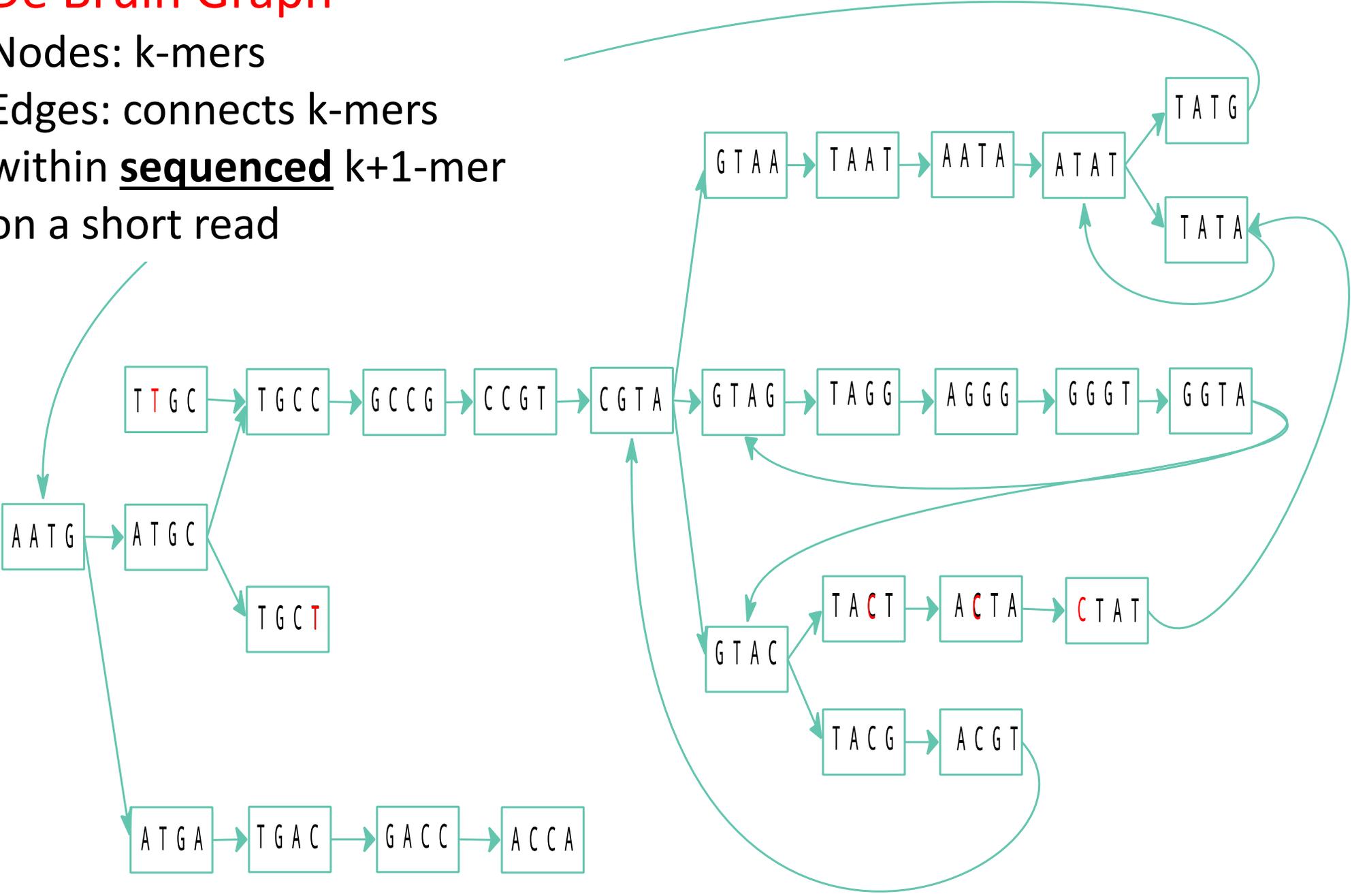


- The city of **Königsberg in Prussia** (now Kaliningrad, Russia) has
 - **two large islands**: Kneiphof and Lomse
 - **two mainland portions**: Altstadt and Vorstadt
 - Connected by **seven bridges**
- **Seven Bridges of Königsberg problem** formulated by **Leonhard Euler** in **1736**
 - Is it possible to make a **walk** (not necessarily a loop) walking **each bridge exactly once**?
 - **Euler** (why the base of natural logarithm is **e**): Basel → Berlin → St. Petersburg
- Euler **proved** that **such walk is impossible**
 - Graph has **three nodes** with **odd degrees**: 3,3,and 5
 - **Only two nodes with odd degrees are allowed**: start and finish of the walk
 - This laid the **foundations of graph theory** and prefigured the idea of **topology**

De Bruin Graph

Nodes: k-mers

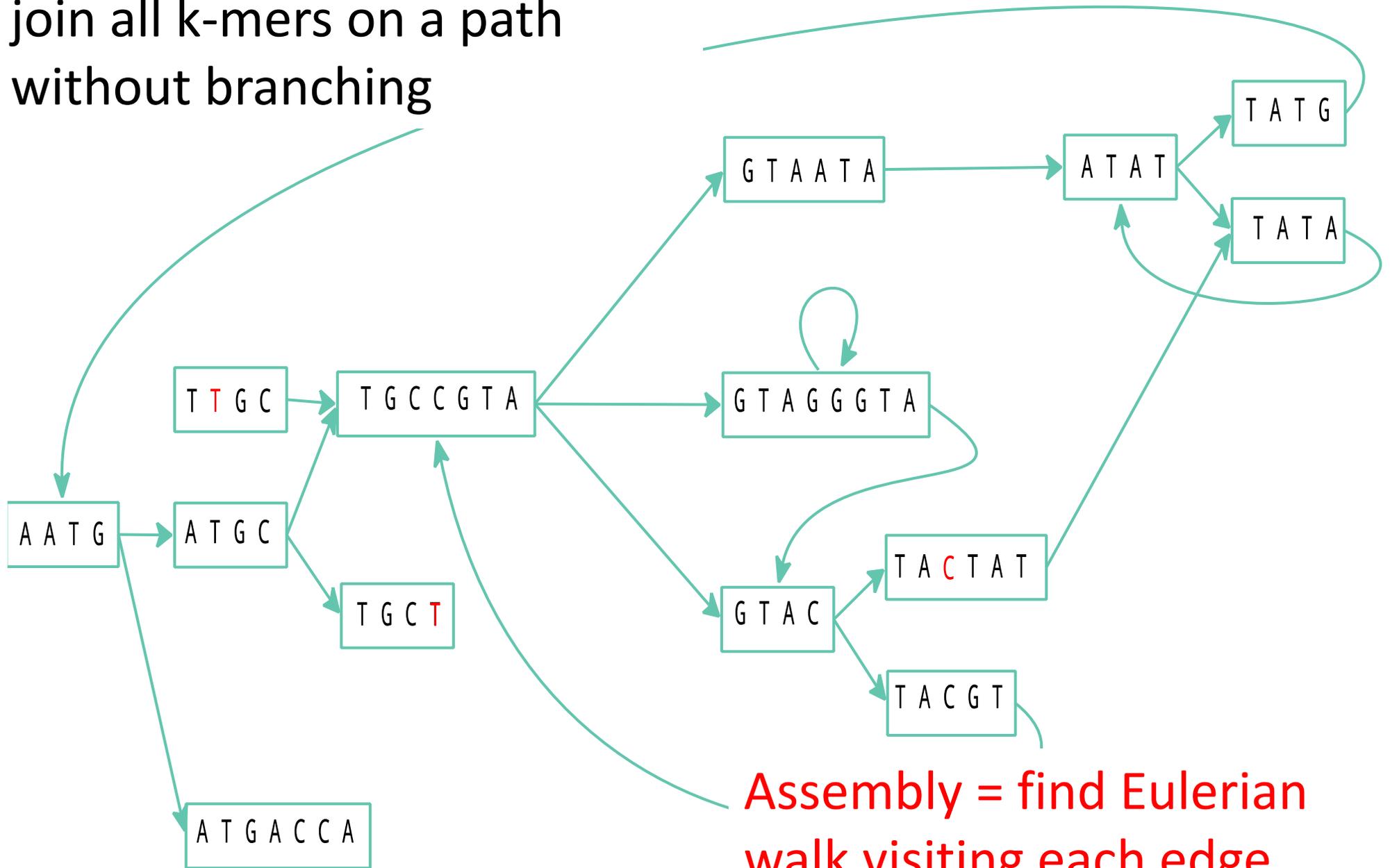
Edges: connects k-mers within sequenced k+1-mer on a short read



Slide by Sorin Istrail, Brown U.

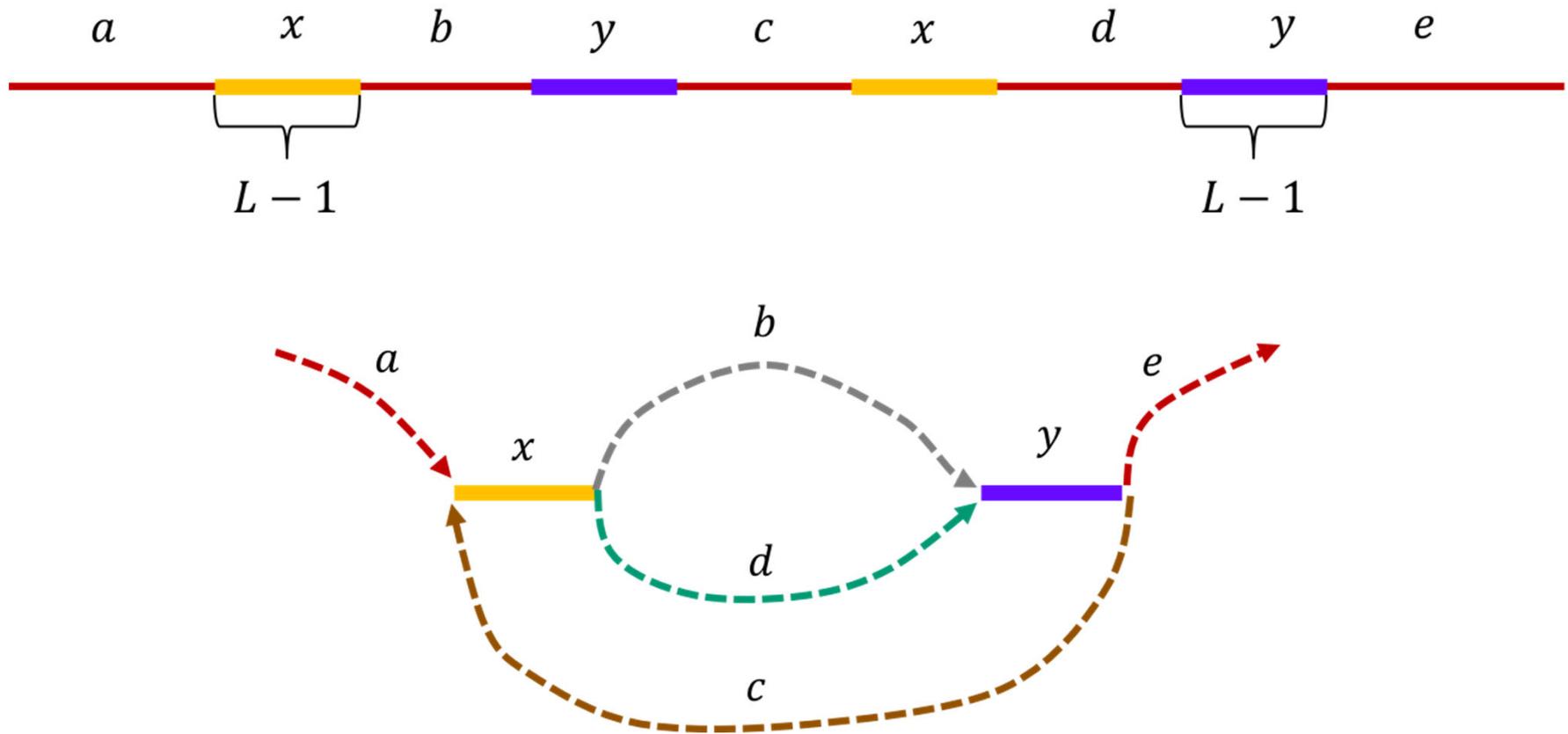
Simplified De Bruin Graph

join all k-mers on a path
without branching



Assembly = find Eulerian walk visiting each edge once

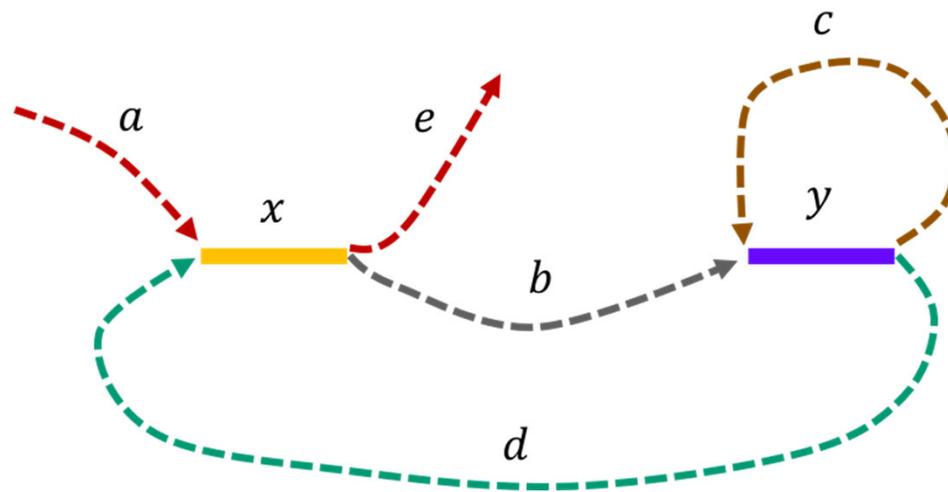
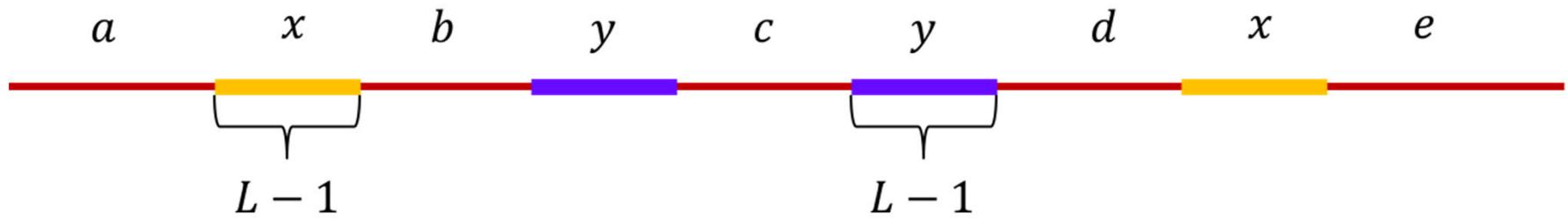
Why interleaved repeats are dangerous?



The two Eulerian paths that are on the graph:
 $a-x-b-y-c-x-d-y-e$ and $a-x-d-y-c-x-b-y-e$

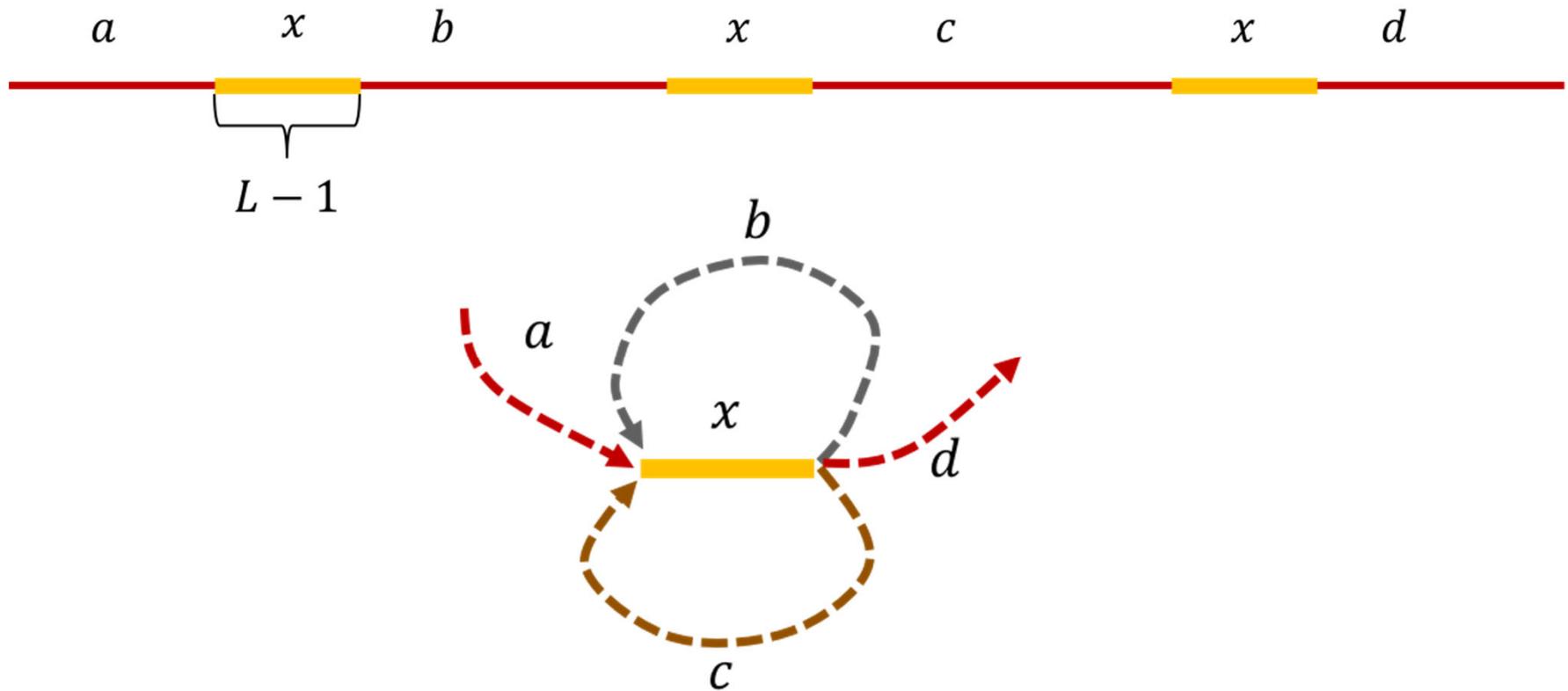
Images from the course EE 372: Data Science for High-Throughput Sequencing.
taught by David Tse at Stanford

Why non-interleaved repeats are safe?



The only Eulerian path is: $a-x-b-y-c-y-d-x-e$

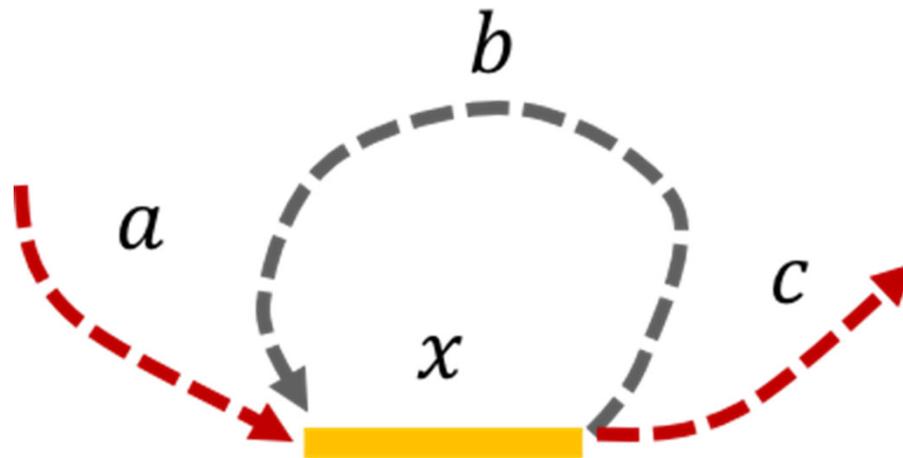
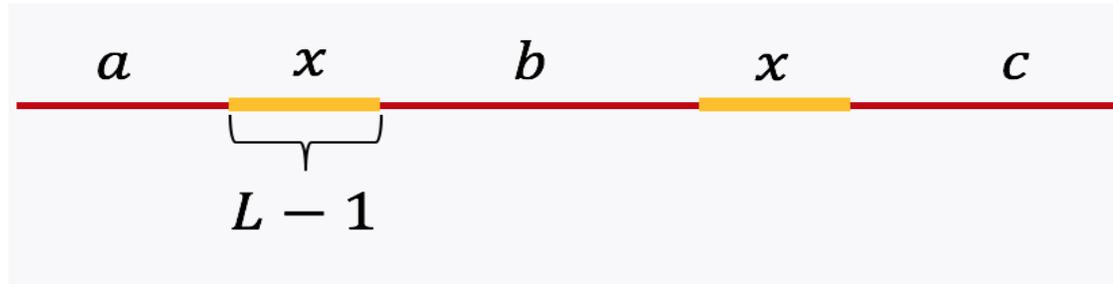
Why triple repeats are dangerous?



The two Eulerian paths that are on the graph:
 $a-x-b-x-c-x-d$ and $a-x-c-x-b-x-d$

Images from the course [EE 372: Data Science for High-Throughput Sequencing](#),
taught by David Tse at Stanford

Why double repeats are safe?



The only Eulerian path is: $a-x-b-x-c$

Pavel Pevzner's theorem

- **Theorem [Pevzner 1995]:**

If L , the read length, is strictly greater than $\max(\ell_{\text{interleaved}}, \ell_{\text{triple}})$, then the de Bruijn graph has a unique Eulerian path corresponding to the original genome.



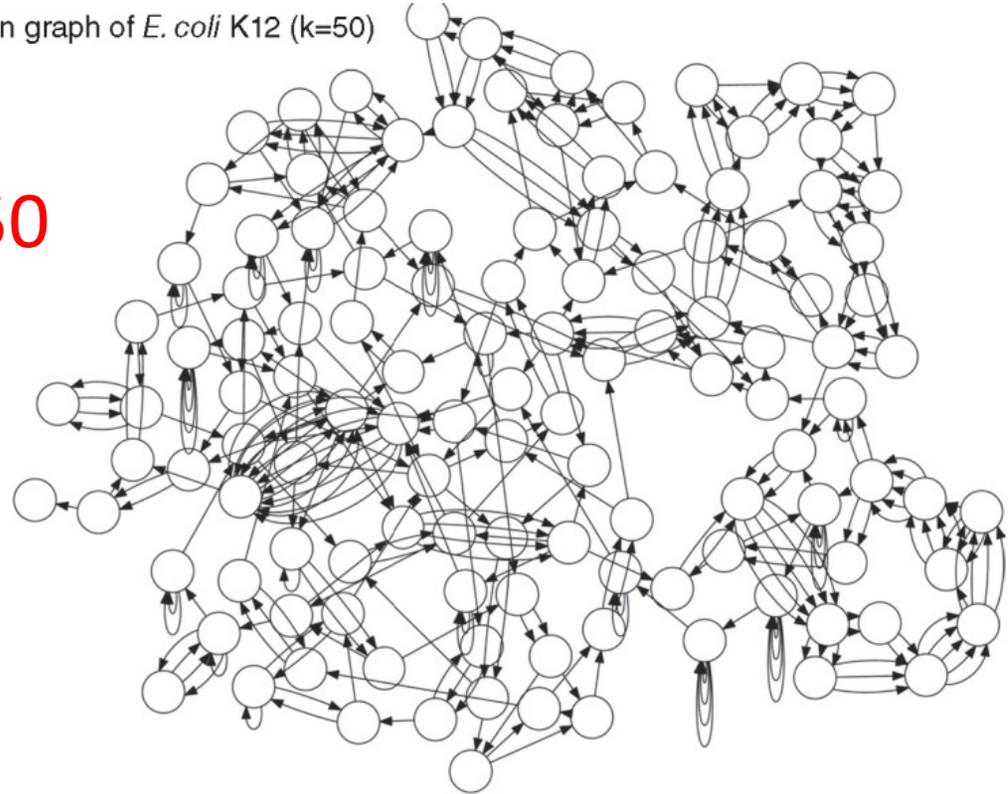
Pavel Pevzner
is the Ronald R. Taylor Chair and
Distinguished Professor of
Computer Science and Engineering
at University of California, San Diego

How to assemble a genome with repeats?

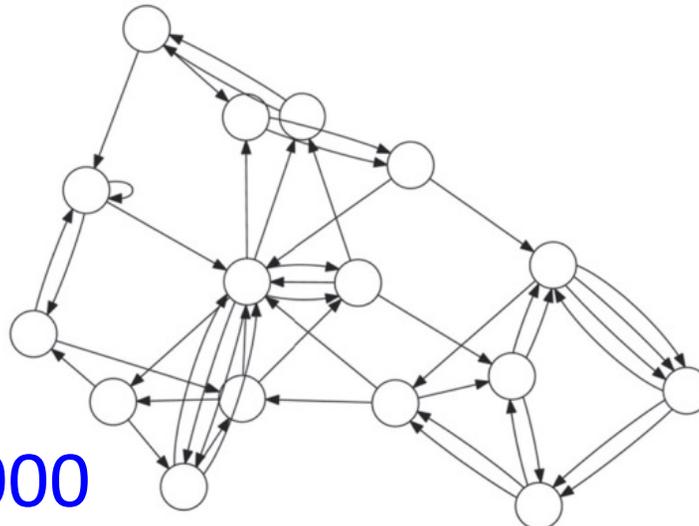
- Answer: use longer reads
- But: cheap sequencing = short reads

(a) de Bruijn graph of *E. coli* K12 ($k=50$)

$k=50$



(b) de Bruijn graph ($k=1,000$)



$k=1000$

(c) de Bruijn graph ($k=5,000$)



$k=5000$

Technology	Read length (bp)
Roche 454	700
Illumina	50–250
SOLiD	50
Ion Torrent	400
Pacific Biosciences	>10,000

A gallery of useful
discrete probability distributions

Geometric Distribution

- A series of **Bernoulli trials** with **probability of success = p** . continued **until the first success**. X is the number of trials.
- Compare to: Binomial distribution has:
 - Fixed number of trials = n . $P(X = x) = C_x^n p^x (1 - p)^{n-x}$
 - Random number of successes = x .
- Geometric distribution has reversed roles:
 - Random number of trials, x
 - Fixed number of successes, in this case 1.
 - Success always comes in the end: so no combinatorial factor C_x^n
 - $P(X=x) = p(1-p)^{x-1}$ where:
 $x-1 = 0, 1, 2, \dots$, the number of failures until the 1st success.
- **NOTE OF CAUTION: Matlab, Mathematica**, and many other sources use x to denote the **number of failures until the first success**. We stick with **Montgomery-Runger notation**

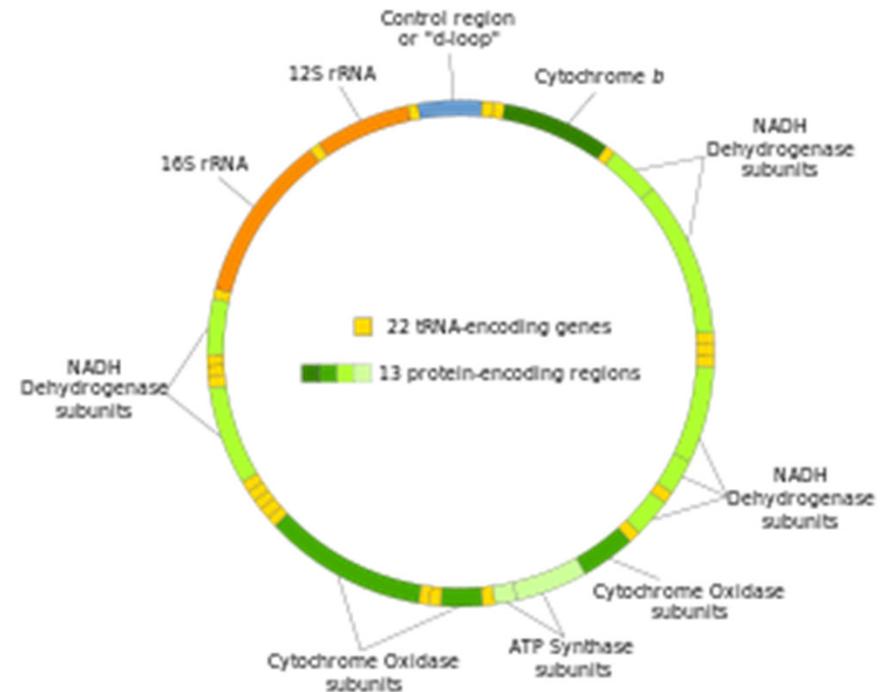
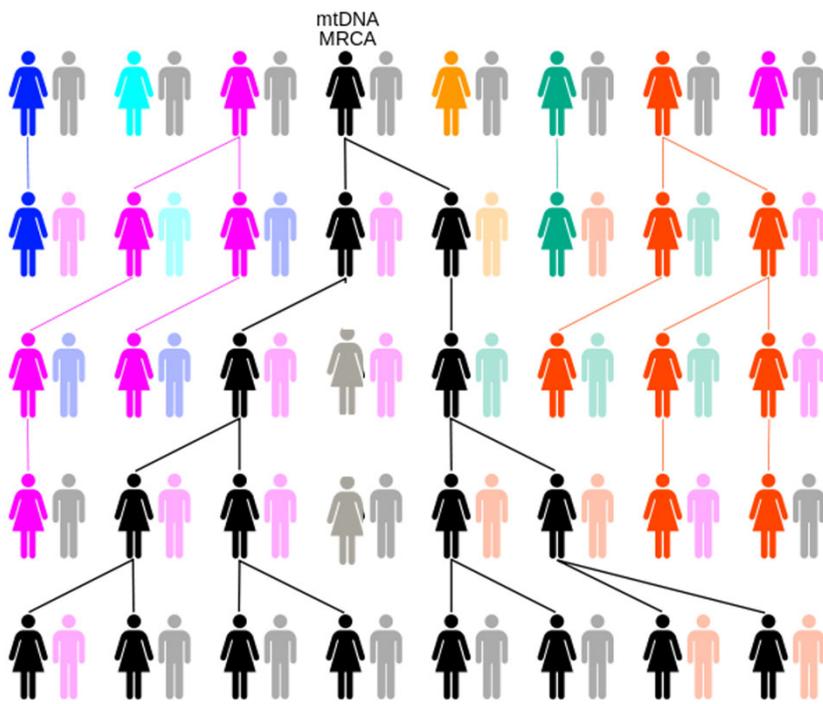
Geometric Mean & Variance

- If X is a geometric random variable (according to Montgomery-Bulmer) with parameter p ,

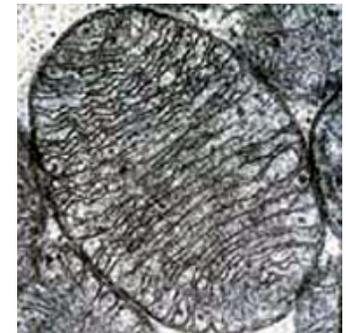
$$\mu = E(X) = \frac{1}{p} \quad \text{and} \quad \sigma^2 = V(X) = \frac{(1-p)}{p^2} \quad (3-10)$$

- For small p the standard deviation $= (1-p)^{0.5}/p \approx$
mean $= 1/p$
- Very different from Binomial and Poisson, where
variance $=$ mean and standard deviation $=$ mean^{1/2}

Geometric distribution in biology

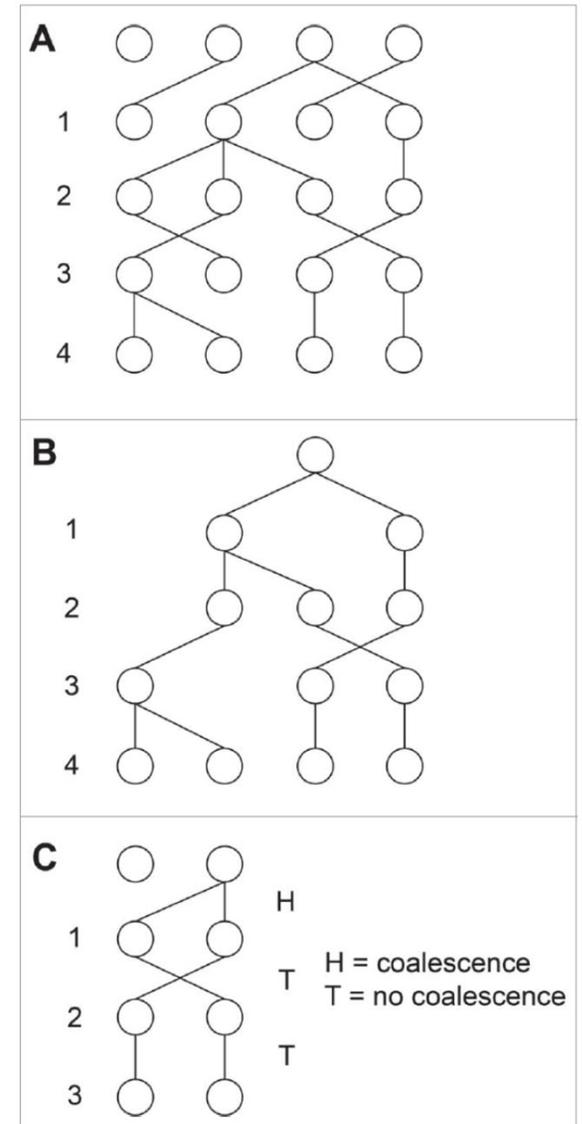


- Each of our cells has mitochondria with 16.5kb of mtDNA **inherited only from our mother**
- Human mtDNA has 37 genes encoding 13 proteins, 22+2 tRNA & rRNA
- Mitochondria appeared 1.5-2 billion years ago as a symbiosis between an alpha-proteobacterium (1000s of genes) and an archaeon (of UIUC's Carl R. Woese fame)
- Since that time most mitochondrial genes were transferred into the nucleus
- Plants also have plastids with genomes related to cyanobacteria



Time to the last common (maternal) ancestor follows geometric distribution

- **Constant population** of N women
- **Random number** of (female) **offsprings**. Average is 1 (but can be 0 or 2)
- **Randomly pick two women**.
Question: how many **generations T** since their **last maternal ancestor**?
- T is a random variable What is its PMF: **$P(T=t)$** ?
Answer: $P(T=t)$ follows a **geometric distribution**
- Do these two women have **the same mother**? Yes: **“success”** in finding their last common ancestor (**$p=1/N$**). **$P(T=1)=1/N$** .
- No? **“failure”** (**$1-p=1-1/N$**). Go to their mothers and repeat the same question.
- **$P(T=t)=(1-1/N)^{t-1}(1/N) \approx (1/N) \exp(-(t-1)/N)$**
- **t** can be inferred from **the density of differences on mtDNA $=2\mu t$**



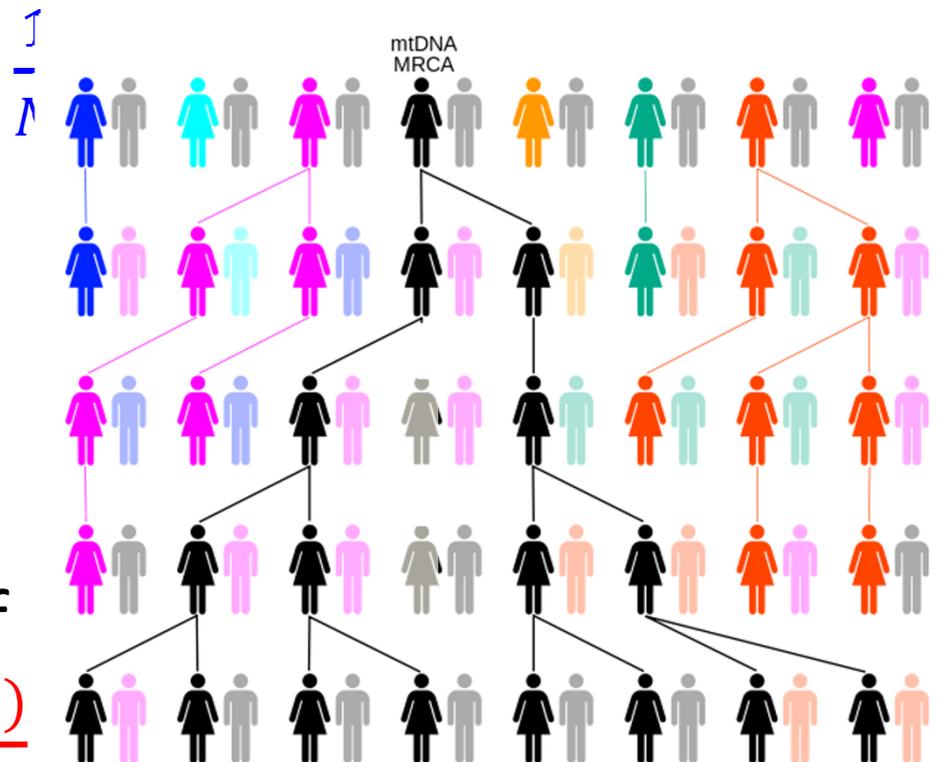
Most Recent Common Ancestor (MRCA)

- Consider N women living today. Let $A(t)$ be the number of maternal ancestors t generations before present connected to all modern women
- At $t=0$, $A(0) = N$. Any of $\frac{A(0)(A(0)-1)}{2}$ pairs of woman can coalesce with probability $\frac{1}{N}$

- We expect $\frac{1}{N} \frac{A(0)(A(0)-1)}{2}$ pairs to coalesce

$$A(1) = A(0) - \frac{1}{N} \frac{A(0)(A(0)-1)}{2}$$

- Now the expected number of coalescing pairs $\frac{1}{N} \frac{A(1)(A(1)-1)}{2}$



Most Recent Common Ancestor (MRCA)

$$\frac{dA(t)}{dt} = - \overbrace{\frac{1}{N}}^{\text{probability of coalescence}} \times \overbrace{\frac{A(t)(A(t)-1)}{2}}^{\text{number of pairs}} \approx \frac{A(t)^2}{2N}$$

Solution: $A(t) = \frac{C}{t+2}$

Let's check: $\frac{dA}{dt} = -\frac{C}{(t+2)^2} = \frac{A^2}{C} \rightarrow C = 2N$

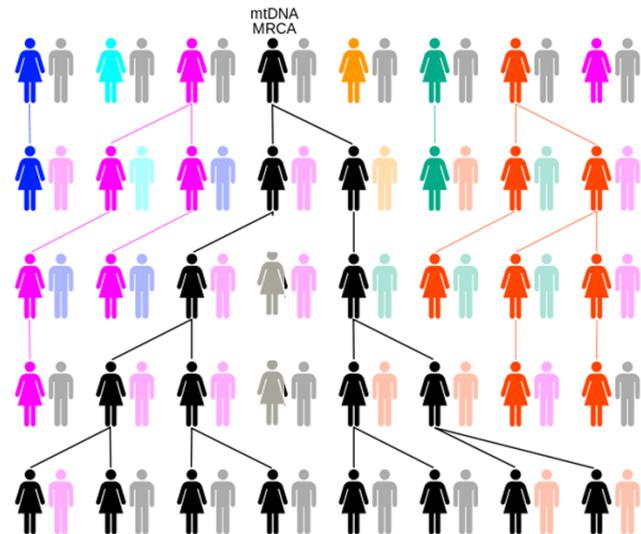
$$T_{MRCA} = 2(N - 1) \approx 2N$$

Most Recent Common Ancestor (MRCA)

- Start with N individuals. Unit of time is N generations (time for one pair to merge) since $E(T) = \sum_{t=1}^{\infty} t \cdot (1/N) \exp(-t/N) = N$
- Any of $\frac{N(N-1)}{2}$ pairs can merge first. The average time for the first pair to merge is $\frac{2}{N(N-1)}$

- After the merger $N \rightarrow N - 1$,

- So, the time until the next merger is longer $\frac{2}{(N-1)(N-2)}$



Most Recent Common Ancestor (MRCA)

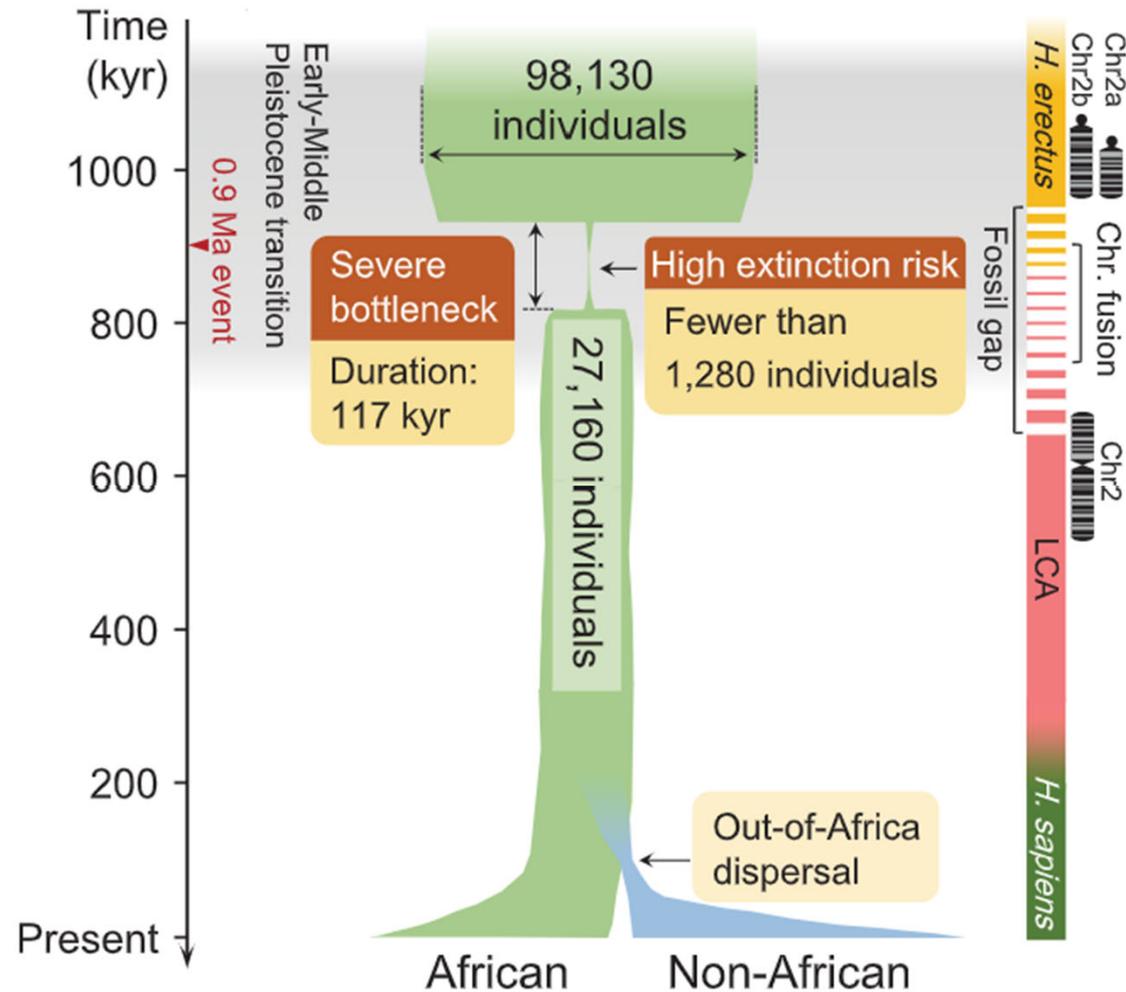
Total time until the MRCA

$$T_{MRCA} = N \cdot \sum_{k=2}^N \frac{2}{k(k-1)}$$

$$= 2N \sum_{k=2}^N \left(\frac{1}{k-1} - \frac{1}{k} \right) = 2N \left(1 - \frac{1}{N} \right) \approx 2N$$

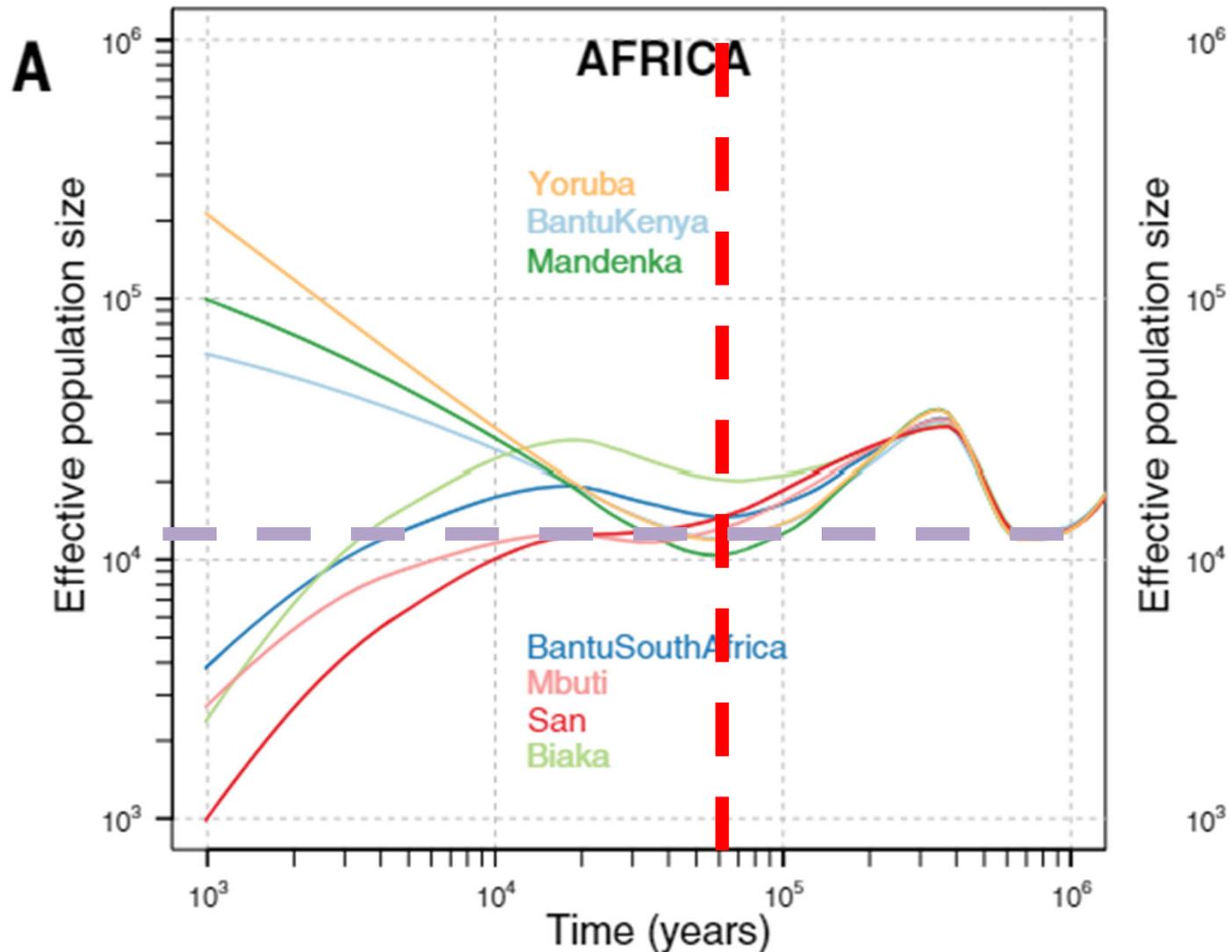
- There are about $N = 8 \times 10^9 / 2 = 4 \times 10^9$ women living today
- **M**ost **R**ecent maternal **C**ommon **A**ncestor (**MRCA**) of all people living today lived $T_{MRCA} = 2N$ generations ago
- $T_{MRCA} = 2 \cdot 4 \times 10^9$ generations
- If the generation time 20 years it is 160 billion years > **10 times the time since the Big Bang.**
- Something is wrong here!

Hot off the press: human ancestors almost got extinct about 1M years ago



Hu W, et al. Science. 2023;381: 979–984

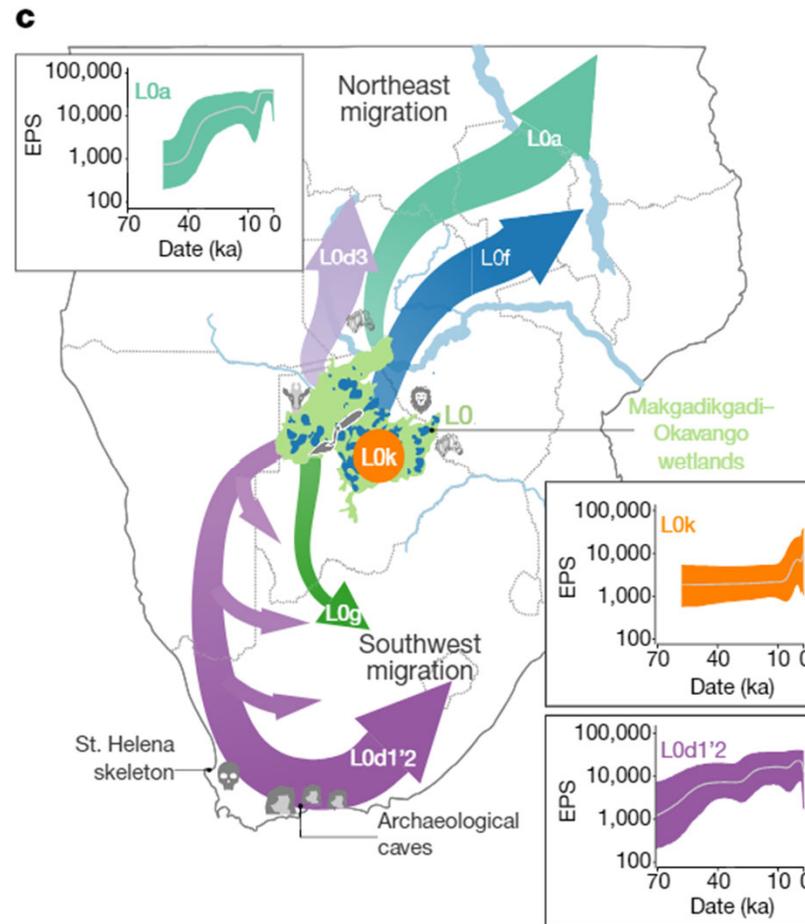
Bottleneck for human population in Africa around 10,000 individuals reached around 60,000 years ago



From ~1000 modern genomes: Bergström A, et al. Science. 2020;367

- Population is **not constant** and for a long time was very low
- Change N to the “effective” size N_e reached during **the bottleneck**
- Current thinking is that for all of humankind $N_e \sim 10,000$ people
- **Mito Eve lived in Africa** $\sim 2 * (N_e/2) * 20$ years = $10,000 * 20$ years = **200,000 years ago**

“Mitochondrial Eve” lived in Africa



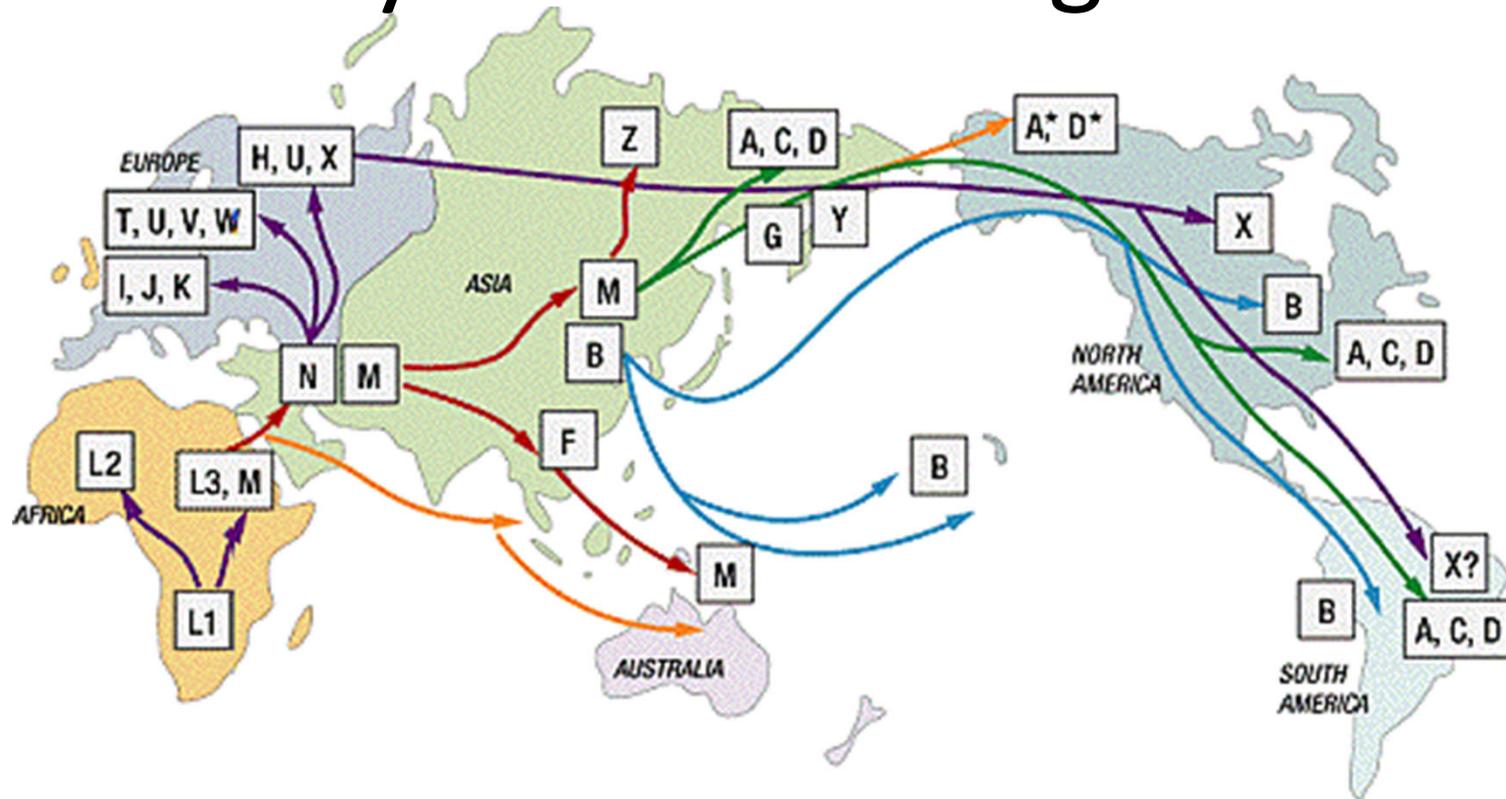
“Mitochondrial Eve” lived in Makgadikgadi–Okavango paleo-wetland of southern Africa ~200,000 years ago (between 165,000 and 240,000 years ago)

Chan EKF, et al. Nature. 2019; 575: 185–189.

Okavango Delta now



Modern mitochondrial DNA contains history of human migrations

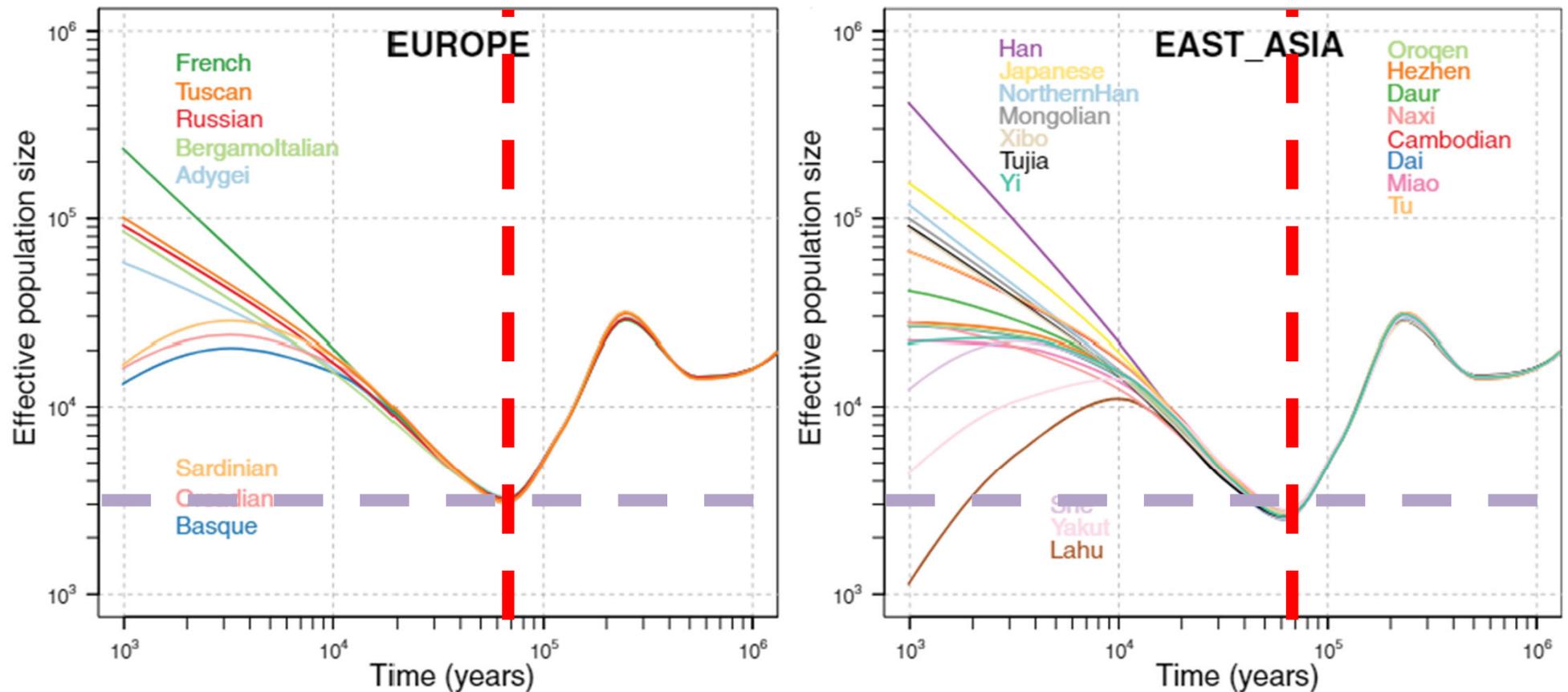


EXPANSION TIMES (years ago)	
Africa	120,000 - 150,000
Out of Africa	55,000 - 75,000
Asia	40,000 - 70,000
Australia/PNG	40,000 - 60,000
Europe	35,000 - 50,000
Americas	15,000 - 35,000
Na-Dene/Esk/Aleuts	8,000 - 10,000



Poznik GD, et al (Carlos Bustamante lab in Stanford), Science **341**: 562 (August 2013).

- Bottleneck for human population in Europe and Asia was around 3,000 individuals reached around 70,000 years ago
- Non-African Eve lived ~60,000 years ago



From ~1000 modern genomes: Bergström A, et al. Science. 2020;367

What about men?

- Y-chromosome is transferred from father to son
- Like mitochondria it can be used to trace ancestry of all men to the “Y-chromosome Adam”
- Where did “Adam” live? Did he meet the “mitochondrial Eve”?

Y-chromosomal Adam also lived in Africa

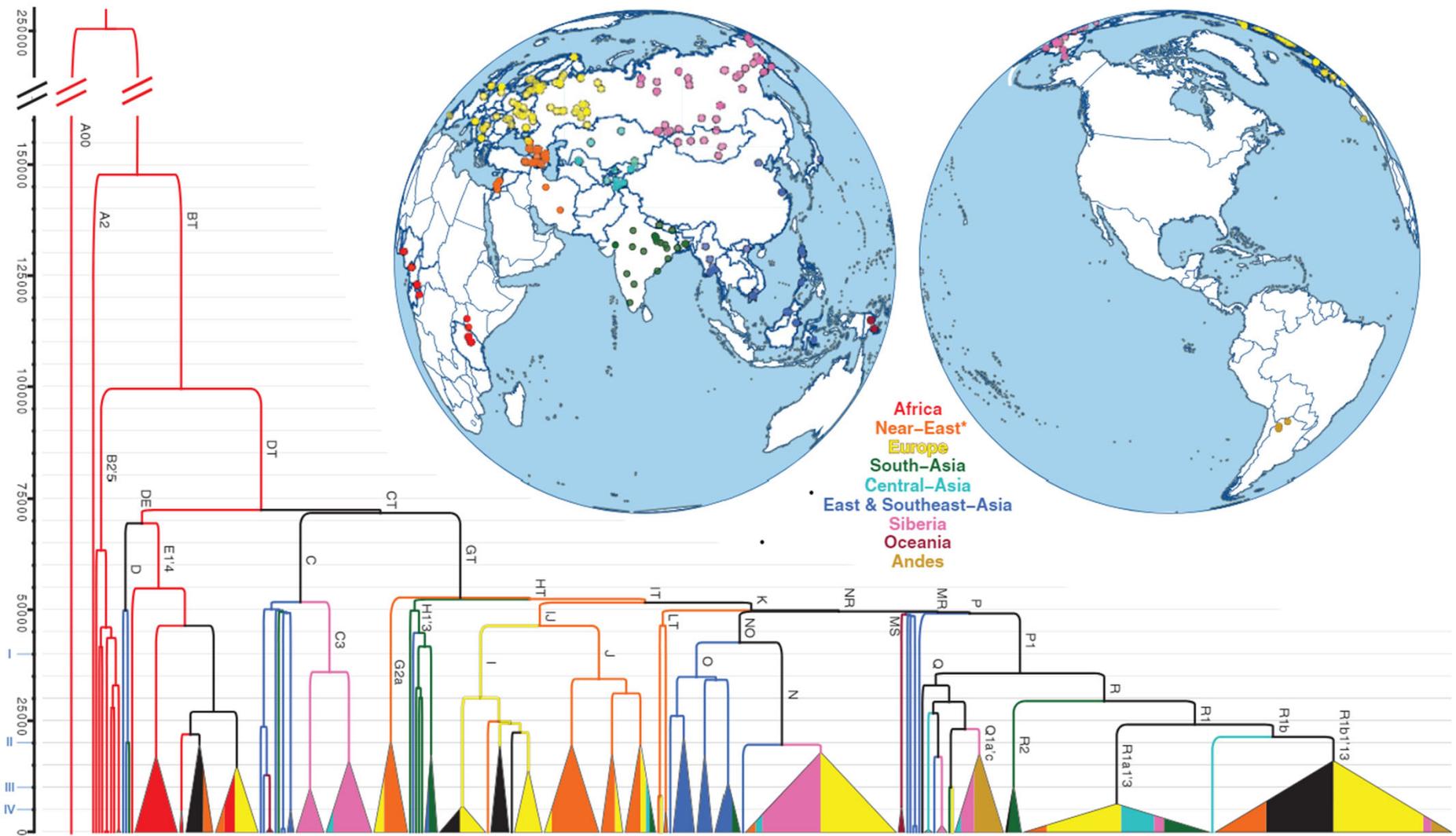
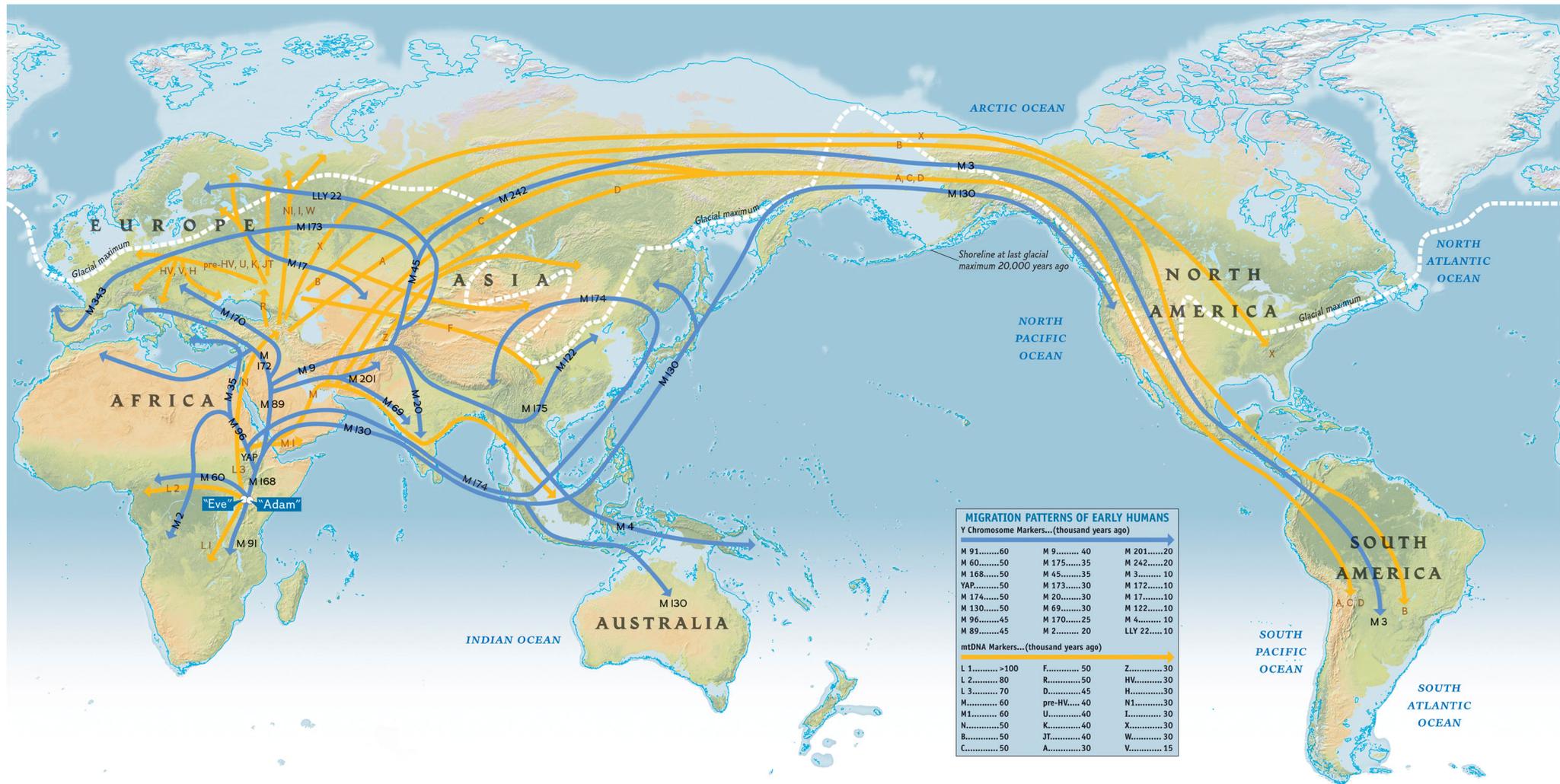


Figure 1. The phylogenetic tree of 456 whole Y chromosome sequences and a map of sampling locations. The phylogenetic tree is reconstructed using BEAST. Clades coalescing within 10% of the overall depth of the tree have been collapsed. Only main haplogroup labels are shown (details are provided in Supplemental Information 6). Colors indicate geographic origin of samples (Supplemental Table S1), and fill proportions of the collapsed clades represent the proportion of samples from a given region. Asterisk (*) marks the inclusion of samples from Caucasus area. Personal Genomes Project (<http://www.personalgenomes.org>) samples of unknown and mixed geographic/ethnic origin are shown in black. The proposed structure of Y chromosome haplogroup naming (Supplemental Table S5) is given in Roman numbers on the y-axis.

Karmin M, Saag L, Vicente M, Sayres MAW, Järve M, Talas UG, et al. *Genome Res.* 2015;25: 459–466.

“Adam” and “Eve” both lived in Africa



- “Mitochondrial Eve” lived in Africa between 100,000 and 240,000 years ago
- “Y-chromosome Adam” also lived in Africa between 120,000 and 160,000 years ago
- Poznik GD, et al (Carlos Bustamante lab in Stanford), *Science* **341**: 562 (August 2013).

Mitochondrial Eve (maternally transmitted ancestry)
Y-chromosome Adam (paternally transmitted ancestry)
lived ~200,000 years ago.

When lived the latest common ancestor shared by all of us based on nuclear DNA?

- A. 1 million years ago
- B. 200,000 years ago
- C. 3400 years ago
- D. 1320 years ago
- E. Yesterday, I really have no clue

Get your i-clickers

Mitochondrial Eve (maternally transmitted ancestry)
Y-chromosome Adam (paternally transmitted ancestry)
lived ~200,000 years ago.

When lived the latest common ancestor shared by all of us based on nuclear DNA?

- A. 1 million years ago
- B. 200,000 years ago
- C. 3400 years ago**
- D. 1320 years ago
- E. Yesterday, I really have no clue

Get your i-clickers

Last common ancestor in nuclear (non Y-chr) DNA is another matter

- Unlike Mito or Y-chromosome, **nuclear DNA gets mixed with every generation**
 - Each of us gets 1/2 of nuclear DNA from the father and 1/2 from the mother
 - Each of us has 2 parents, 4 grandparents, 8 great-grand parents ...
- If one assumes:
 - Well-mixed marriages (not true: mostly local marriages)
 - Constant size population (not true: much smaller in the past)
 - In 33 generations the number of ancestors:
 $2^{33} = 8 \text{ billion} = 8 \text{ billion people living today}$
- Every pair of us living today should have at least one shared ancestor who lived
 - 33 generations * 20 years/generation = 660 years ago ~1360 AD
- **Assuming $T_{MRCA} = 2 T_{\text{average pairwise ancestor}} = 2 \cdot 660 \text{ years ago} = 700 \text{ AD}$**

Corrected for (mostly) local marriages and rare migrations

and rare migrations

Modelling the recent common ancestry of all living humans

Douglas L. T. Rohde¹, Steve Olson² & Joseph T. Chang³

¹Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

²7609 Seabrook Road, Bethesda, Maryland 20817, USA

³Department of Statistics, Yale University, New Haven, Connecticut 06520, USA

With 5% of individuals migrating out of their home town, 0.05% migrating out of their home country, and 95% of port users born in the country from which the port emanates, the simulations produce a mean **MRCA date of 1,415 BC** and a mean **IA date of 5,353 BC**.

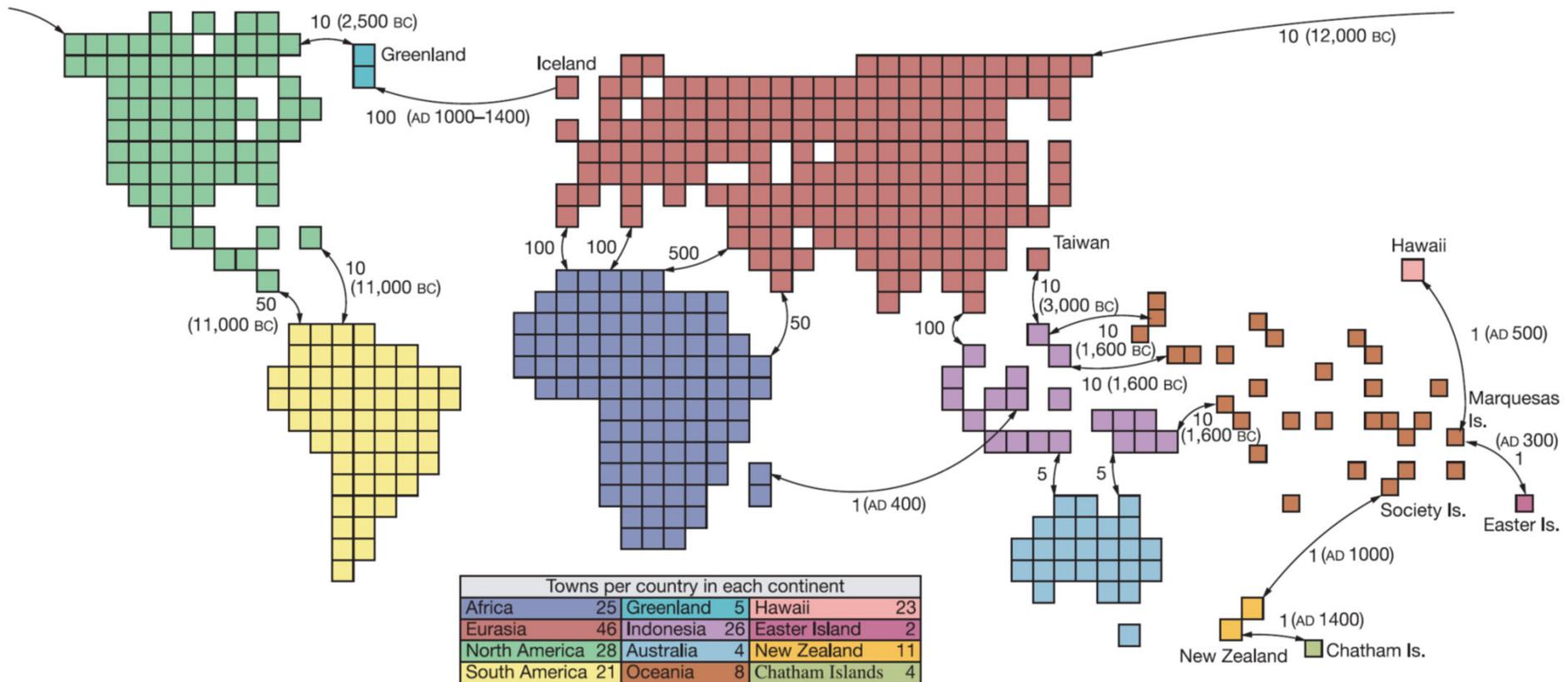


Figure 2 Geography and migration routes of the simulated model. Arrows denote ports and the adjacent numbers are their steady migration rates, in individuals per generation. If

given, the date in parentheses indicates when the port opens. Upon opening, there is usually a first-wave migration burst at a higher rate, lasting one generation.